

DiscoFuzzer: Discontinuity-based Vulnerability Detector for Robotic Systems

Sean Rivera ¹, Antonio Ken Iannillo ², and Radu State ²

¹University of Luxembourg

²Affiliation not available

October 30, 2023

Abstract

We demonstrate a new fuzzing tool for the Robotic Operating System (ROS), which exploits the physical nature of robotic systems to detect a novel class of bugs.

DiscoFuzzer: Discontinuity-based Vulnerability Detector for Robotic Systems

Sean Rivera
University of Luxembourg
sean.rivera@uni.lu

Antonio Ken Iannillo
University of Luxembourg
antonio.iannillo@uni.lu

Radu State
University of Luxembourg
radu.state@uni.lu

ABSTRACT

Robotic systems continue their diffusion into society, accomplishing a myriad of physical tasks on our behalves. However, their safety-critical nature implies the vitality of testing their robustness. In this paper, we propose a novel fuzzing methodology that exploits the continuity of the physical world to automatically explore the input space and detect malfunctions in robotic software modules. The analysis of outputs includes the computation of the first and second derivative functions that can unveil anomalies in the behavior of the software. We implemented this methodology in *DiscoFuzzer*, the discontinuity-based fuzzer for ROS, and evaluate three different sampling approaches based on Monte Carlo, Chebyshev, and Spline methods. *DiscoFuzzer* detected 85 distinct vulnerabilities: 77 of the 89 previously known vulnerabilities, and 8 novel vulnerabilities previously undetected. The discontinuity analysis of *DiscoFuzzer* detected 41 more unique vulnerabilities than crash detection alone. Furthermore, we determined that the implemented sampling approaches were statistically equal in finding vulnerabilities, but each of them discovered vulnerabilities the others could not. The chebfun sampling is found to be the fastest in finding vulnerabilities. We determined that *DiscoFuzzer* provides a unique fuzzing solution for ROS systems to detect a wide variety of security vulnerabilities with high precision.

1 INTRODUCTION

Robotic systems are proliferating in our society due to their capacity to carry out physical tasks on behalf of human beings. Current applications include, but are not limited to, military, industrial, agricultural, and domestic robots[6]. In this paper, we consider robotic systems as a subset of cyber-physical systems (CPS). While CPS are defined as “physical and engineered systems whose operations are monitored, controlled, coordinated, and integrated by a computing and communicating core”[32], a robot is a “actuated mechanism programmable in two or more axes with a degree of autonomy, moving within its environment, to perform intended tasks”[14]. The standard ISO8373 distinguishes robots further in industrial and service robots according to their tasks.

However, in the deployments of robots it is vital to consider their safety-critical nature. Indeed, a faulty robot can irreversibly damage the physical environment in which it is operating, including being harmful to human beings. Thus, testing the robustness of robotic systems is crucial to ensure safety.

Fuzzing has become a central tenant of computer security research over the past 20 years. Fuzzing is the repeated execution of the program using inputs sampled from the input space [22]. This repeated execution is performed automatically with the sampling of inputs adhering to different heuristics. Usually, fuzzers guided by

code-coverage [9, 16, 40, 44, 55] or by symbolic engines dominate the scene of the security community [7, 46, 57, 59]. Furthermore, most current fuzzers solely focus on traditional software crashes or well-known vulnerable code behaviors with the help of sanitizers [20, 21, 41].

However, robotic systems have to face different threats that current fuzzers do not adequately address. In this paper, we consider external security threats that consist of an attacker gaining control of specific modules of a robotic system or interfering with the robot’s sensors input data[10, 11, 19, 34]. An attacker can lead a robotic system to disregard safety measures, behave in an unsafe way, or crash.

We propose a **novel fuzzing methodology that exploits the continuity of the physical world to automatically explore the input space and detect malfunctions in robotic software modules**. Our solution improves over conventional software testing because it directly addresses failure states, *i.e.*, unsafe behaviour, that other fuzzers do not have the ability to address, and provides a mathematical approach to provide dynamic solutions. The fuzzer operates by performing a single test against a target and constructing a model of the targets’ response behavior, including discontinuity analysis. We show how a comparison of these tests on several inputs can detect anomalies that are usually only unveiled by extensive manual testing.

We implemented *DiscoFuzzer*, the discontinuity-based fuzzer for ROS[31]. We focus on the Robotic Operating System (ROS) because it is an open-source project with a large active community. ROS is a software framework that provides a structured communication layer and libraries to develop robotic applications. Manufacturing and other industries[35][37] use ROS for their robots, while both Windows and Amazon have shown interest in supporting it[54][42]. With ROS systems projected to make up the majority of robotic systems within the next five years, a focus on security is needed [28]. To the best of our knowledge, *DiscoFuzzer* is the first fuzzer applied to ROS. However, Santos *et al.* proposed a property-based testing approach for ROS. Property-based testing is very similar to fuzz testing, but the latter does not require familiarity with the target to specify *properties* and *shapes*[39].

Our main contributions in this paper are:

- A novel fuzz testing methodology to examine robotic software systems based on numerical analysis and discontinuity localization;
- A formalized benchmark for ROS consisting of 14 open-source packages and 89 previously known vulnerabilities (plus eight new vulnerabilities discovered by *DiscoFuzzer*);
- The design, implementation, and evaluation of *DiscoFuzzer*, its discontinuity-based detection mechanism, and three different sampling approaches.

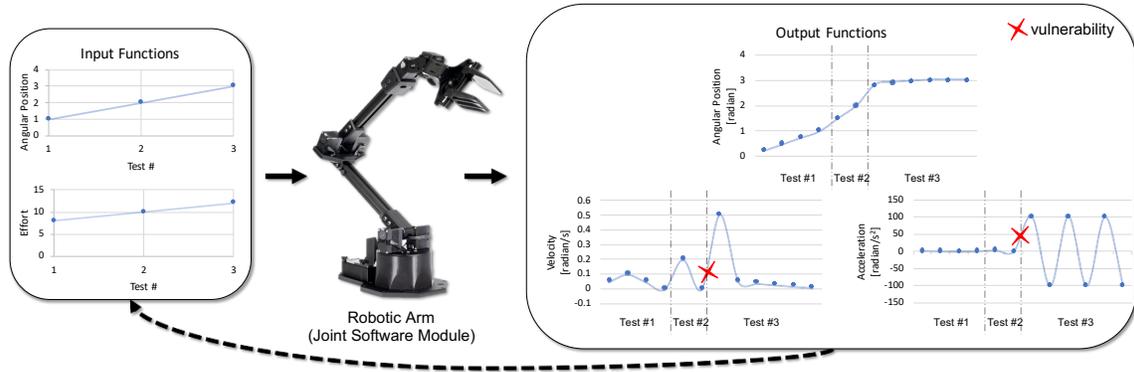


Figure 1: Example of discontinuity-based vulnerability detection for the software module of a robotic arm’s joint.

We demonstrated that *DiscoFuzzer*’s discontinuity analysis could detect more than three times the vulnerabilities that the crash detection mechanism finds. Both *DiscoFuzzer* and the used benchmarks will be publicly released on GitHub for the reproducibility of the experiments and to aid further research in automated testing of ROS packages.

2 DISCONTINUITY-BASED FUZZ TESTING

Let us consider an industrial robotic arm with a single joint. The joint software module receives as input the desired angular position and effort (power) to provide to the joint motor to perform the movement. The output is an array describing the actual change of angular position, velocity, and acceleration in time as performed by the joint. Figure 1 shows an example of how the proposed methodology works. The discontinuity fuzzer begins by first choosing a range of continuous input values. In this case, the range consists of three values for the angular position and three values for the effort. They constitute three fuzz test cases that are provided one by one to the joint software module of the robotic arm. Each defined test case wants to spin the joint one radian clockwise but with increasing effort, *i.e.*, the power to provide to the motor to perform the movement. After each test, we collect the output messages and concatenate their values to create the output functions. The fuzzer then computes first and second derivative functions to analyze discontinuities and unveil anomalies in the behavior of the software. Indeed, the output functions of velocity and acceleration contain a significant change of the slope. The fuzzer detects an anomaly that consists of the uncontrolled rapid rotation of the joint, which could lead to a burn out of the motor because the software did not check the physical constraint. The methodology reports the properties of the output function to help guide input test-case generation.

This example illustrates the functionality of discontinuity fuzzing. Initially, our methodology identifies all of the target’s inputs and outputs. For each input, it creates a sample generator to provide intervals, *i.e.*, ranges of continuous values, depending on the input types (including arrays). Regarding floating points, the sample generator also needs to define the distance between values in the interval (*i.e.*, the resolution). The methodology changes one input at a time while setting the others to default values. Thus, a fuzzer should perform the following actions:

- Choose an interval: The fuzzer samples a range of continuous values and generates this iteration’s test cases (§2.1).
- Test values in the interval: The fuzzer performs one test at the time against the target and saves the outputs.
- Analyze the output distribution: The fuzzer analyzes all of the outputs of this iteration to detect anomalies (§2.2).

The fuzzer iterates on these actions until a termination criterion is met, such as a time limit.

In the following subsection, we present some approaches we adapted to implement the above methodology.

2.1 Sampling approaches

A sampling generator consists of an iterative process to choose the next group of samples to test. We defined three different sampling approaches based on three numerical sampling functions known in the state-of-the-art literature.

First, we used a simple randomized Monte Carlo analysis[43] that chooses every single value at random. Second, we perform a Chebyshev approximation[13] of the target’s input/output function and use the discontinuity and root-finding of the Chebyshev approximation to look for potential areas of interest. Third, we construct a model of the behavior function using a Cubic Hermite Spline Interpolation. We chose these three functions to cover different spectra of potential erroneous behavior.

The sampling approach is different between arrays and individual values, and we expand on them separately below.

Monte Carlo Sample Generator. Monte Carlo is a method based on repeated random sampling. To conduct Monte Carlo sampling for individual values, the sample generator chooses a random value to be the central point for evaluation, and two values are sampled in either direction for derivation analysis, for a total of five values per iteration. For arrays, the fuzzer randomly chooses between (1) creating the array repeating a single sample point, and (2) sampling a central point and creating the array with increments of the central point. Within this paper, we shall refer to Monte Carlo sampling as `monte_carlo`.

Chebyshev Sample Generator. Chebyshev is a numerical method

that generates increasingly complex polynomials (Chebyshev polynomials) for the interpolation of a dataset. Since the initial introduction of the `chebfun` package in Matlab in 2005, Chebyshev polynomials have gained broad interest in the field of numerical approximation research as a way to model an arbitrary function. One of the vital functions that Chebyshev polynomials provide is the efficient location of roots and inflection points, computed at each approximation. The Chebyshev sample generator initially returns a list of interpolation inputs. Once the target runs the interpolation inputs, the new list of interpolation inputs is generated based on the target outputs. Indeed, this approach recursively determines which points would be most optimal to sample, given the current function approximation. Within this paper, we shall refer to Chebyshev sampling as `chebfun`.

Spline Sample Generator. Similar to Chebyshev, Spline is a numerical method that interpolates the function through a piecewise polynomial called a spline. Spline sampling begins with the generation of two clusters of random values, in the same way as `monte_carlo` does. Then, it generates a cluster at the midpoint between these two. If the output from the third selected cluster is similar to the spline interpolated output of the first two clusters, within a user-specified tolerance, then a spline is drawn between the first two clusters. Thus, it randomly generates a new cluster. Otherwise, it generates a new cluster within the bounds set by the first two generated clusters, and it performs a new interpolation. For arrays, it follows the same behavior as `monte_carlo`, with the caveat that the methods for filling arrays are treated as two distinct sources of input for interpolation and generation separately. Within this paper, we shall refer to Spline sampling as `spline`.

Structurally, the `monte_carlo` simulation has the lowest processing overhead and the fastest turnaround, allowing it to cover a larger number of potential inputs. However, the `chebfun` analysis provides a more precise analysis. The `spline` approach acts as a hybrid of the two, with moderate precision and processing overhead.

2.2 Discontinuity Analysis approach

We believe discontinuities in the output functions are the most significant source of issues as the difference in state course changes incredibly quickly as the system runs. The discontinuity analysis approach first looks for discontinuities in the output. It focuses on sharp changes in the slope of the output functions: large values in the first derivative indicate a significant change, and large values in the second derivative indicate a significant effect. We believe that these changes in the derivatives are more likely to be indicative of errors than any other behaviors, as it is less likely for a numerical output with a smaller average derivative change to be the source of vulnerability.

However, there could still be a chance that a minor change in derivative will create a significant effect. As cyber-physical systems tend to maintain a consistent internal state about the world which they use to make decisions, the approach also considers the amount of influence past inputs have on current outputs. This issue is more obvious to recognize with the input of large floating-point values: these values continue to affect the output regardless of the

new provided value due to the mathematical rounding properties. Once an input corrupts the internal state, unexpected or malicious behavior becomes more likely in the long-term. Thus, the analysis raises an error when old inputs in the memory buffer start being disproportionately different from new inputs.

Furthermore, we choose to look for anomalies where changes to one input cause changes in multiple outputs. We enhance the analysis by considering how many times a field in the output topic has changed. If it changes at some consistent rate, then there is presumed to be a connection between the published topic and the output field. However, if there is a one-off outlier, it is more likely to be a memory or mathematical calculation error. For example, consider a memory-overflow vulnerability. If an input makes the software write beyond the bounds of an array, it could overwrite the other inputs. Thus, a field could change in the output response message when it would typically not be affected by that type of input.

In summary, our discontinuity analysis focuses on two main patterns: values changing very quickly and values changing unexpectedly.

In both experimental analysis[45] and industrial surveys[53] algorithmic bugs are found to make up between **30 and 50%** of the causes for robotic software failure. Additionally, previous exploratory research found that more than **75%** of all algorithmic bugs caused either a discontinuity in output (>50%) or a system crash(>20%)[45]. The modular nature of most robotic systems means that while there may be discontinuities in the output of a node, all downstream nodes need to adapt to these discontinuities. For example, a user may create a command to have a robot reset its position to account for rotation. In the real world, this may be advantageous to avoid errors after long runs. However, this does lead to a discontinuity in robot position, which may cause bugs to manifest in other components in the robot, potentially leading to robot burn-out or another failure. Even in these cases, where the user may see an immediate advantage, the importance of addressing an anomaly is vital in avoiding potential failures in other components.

3 ROS CHALLENGE

Robotic Operating System (ROS) is a collection of software libraries that enables communication of both (abstracted) hardware and (pure) software components to develop robotic systems. It is predicted that ROS centered systems will make up the majority of robotic systems within the next five years, both in commercial and academic settings[28]. Even among non-ROS systems, many design similarities exist wherein similar methods can be applied [12]. ROS is an open-source project, and it has a vibrant community with thousands of developers and over nine thousand unique packages available¹. Usually, these packages extend ROS core functionalities by implementing commonly used robotic software modules such as hardware drivers, robot models, datatypes, planning, perception, localization, simulation tools, and other algorithms. Of particular mention is the ROS-Industrial[35] consortium: it consists of 78 international industries that chose to extend and adopt ROS for their industrial robots.

¹<https://metrics.ros.org/>

In ROS, independent computing processes called nodes communicate through messages. Messages define clean and consistent interfaces. A **ROS node** is a self-contained process that controls a part of the robot’s operation. A **ROS topic** is an implementation of a channel in a publish-subscribe model. Topics act as a named bus where nodes can join as either a publisher, a subscriber, or both. When a publisher node sends a message over a topic, every subscriber node to that topic receives a copy.

As an example, Figure 2 shows a node and its topics from the Universal Robot package². The `arm_controller_spawner` plans the movement of a joint to reach a determined position and makes the joint perform this movement. The node subscribes to one topic in which a message consists of different components (inputs), e.g., angular position and effort. It publishes a single topic in which a message consists of different components (outputs): the current angular position, the velocity, and the acceleration. We based the example in Figure 2 on a simplification of the Universal Robot package.

The ROS project already includes guidelines for testing[36] by applying existing libraries, i.e., `gtest` for C nodes and `pytest` for Python nodes, for simple unit testing. A tester should design and implement every test in the proper language with test inputs and oracles for every behavior they want to check. Santo *et al.*[39] extended ROS unit tests with property-based testing. Their tool randomly and automatically generates test inputs while checking two simple properties of the node: liveness (a node should not crash) and interface stability (the set of topics should not change). However, adding new properties requires formal reasoning about the target node to formulate its specific correct behavior.

Thus, we recognize the following challenge: designing a ROS fuzzer to minimize human effort, explore the input space automatically and efficiently, and detect any deviation from correct behavior.

²http://wiki.ros.org/action/show/universal_robots



Figure 2: The Universal Robot’s node and topics

4 DISCOFUZZER

DiscoFuzzer implements the presented novel methodology to target ROS nodes. The fuzzer is developed in Python 3.6, and targets ROS kinetic.

Figure 3 shows an overview of *DiscoFuzzer*. The fuzzer’s input is a user-defined configuration file that includes target information. It outlines the type of messages used in the topics published or subscribed to by the target node. It defines the input topic for the sampling generators. The configuration file also contains how to execute the target node, i.e., flags, and arguments for the command-line. An orchestrator reads the configuration file and it ① initializes the ROS environment with the target node and *DiscoFuzzer*’s publisher and subscriber. A timer starts to count how long the fuzzer should test the target. Then, the first iteration starts. The orchestrator ② notifies the test case generator of the new iteration. The generator samples a new interval, and it ③ provides one value at a time to the ROS environment through the publisher. Once the subscriber receives a message, it ④ copies this value into the output analyzer. The analyzer ⑤ sends the result of the analyses back to the orchestrator. When anomalies are detected, the orchestrator ⑥ relaunches the node and returns the node to a starting state so that testing can continue. The fuzzer’s output is a set of reports issued during the fuzzing campaign.

4.1 Test case generator and Publisher

The test case generator takes the current state of the system and generates a new set of inputs to test through the publisher. The user can define the *interval size*, i.e., the number of inputs to generate at every iteration. If the sampling generator computes the size algorithmically (e.g., in `chebfun`), it ignores this configuration parameter.

The generator tracks which inputs it has already given to the system to ensure that there are no unnecessarily repeated values. Beyond that, it maintains a collection of every single combination of potential inputs from the topic to publish on. This collection allows us to fuzz individual inputs one at a time, as well as inputs in groups of size $s_2 \dots s_n$, where n is the size of the input. It allows for the identification of potential interactions or discontinuity in situations where the combined two inputs can break the system, even if they could not individually. It focuses on single inputs and implement a power-law drop-off for values fuzzed in combination,

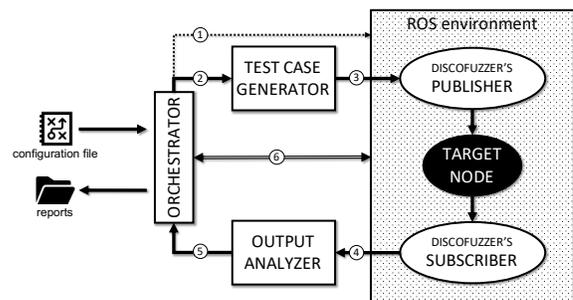


Figure 3: Overview of *DiscoFuzzer*

prioritizing potentially interesting interactions found from the single fields before conducting a random search. For every ten values fuzzed, nine values are single, and one is a combination. For the combinations, 90% are pairing, 10% are combinations greater than two.

The test generators implement the three sampling approaches presented in the methodology.

For the `monte_carlo` sampling approach, the fuzzer chooses the central input at random. The fuzzer stores every point, and it discards any new point that overlaps with the previous range. Additionally, we chose to initialize the fuzzer with some common failure-inducing values such as NaN, infinity, and very large or very small floating points.

Since the `chebfun` approach works by calculating a polynomial and measuring the result, the generator relies on the information provided by the output analyzer (through the orchestrator). For our implementation, we chose to use the `pychebfun` library[47] as it had the most active development and stars on GitHub outside of the Matlab core implementation. As the `chebfun` approach stores only the expansion coefficients we do not require any space-saving optimizations for a longer run.

The `spline` approach consists of a numerical analysis that stores a group of splines containing both measured points and their derivatives, allowing for more accurate interpolation. One of the main benefits of this method is that it allows arbitrary N-dimensional interpolation of functions without extra computation. The `spline` interpolation allows us to perform a more accurate analysis of array fields. For our implementation, we rely on the `scipy` interpolation library [49], choosing our points for the initial array in a random method, just like the `monte_carlo` approach. The `spline` interpolation approach is refined with a mixture of midpoints and random additional points.

For generating an interval of floating points, `DiscoFuzzer` requires the user to choose a *resolution*. The resolution both represents the distances between points in the same interval and indicates the number of decimal digits to consider.

The rate of publishing is limited in two areas: (1) the rate at which the targeted node can respond and (2) the rate at which the subscriber can update. While we can optimize the subscriber and assign new threads to it, there is nothing we can do about the targeted node. As such, the publisher runs on a single thread, and no additional processing capabilities are required.

Users can also choose to set limitations on the `monte_carlo` sampling if they expect that only specific values will produce interesting results. Users can set a *compression factor* and/or place *constraints on combinations*.

In the first case, `DiscoFuzzer` will choose a range at random with size equal to the compression factor with bounds determined by the *resolution*. Then, `DiscoFuzzer` sweeps through the range from bottom to top at the chosen resolution, only saving the center of the range instead of saving each selected random point. Effectively, whenever `DiscoFuzzer` chooses a value for its sampling, it chooses the next value in the specified range instead of a completely random value. Once the range has been fully explored, it will randomly select another range and repeat the sweeping process. This gives users far more space-efficient storage for their simulation at the

cost of far fewer random values being explored over time. As an example, suppose the user sets a range size of 100. When using the *compression factor* feature, after sampling 10,000 points it would only have 100 values saved in memory instead of the 10,000 total values it would have under normal sampling procedures.

The *constraints on combinations* feature allows users to set constraints on generated inputs. These can vary in complexity from simple constraints such as limiting fields to a specific range of values (such as $x < N$), to specific combination rules (such as $y > x; y < N$). Users can specify any numeric constraints they deem necessary for their purposes. This allows users to better guide fuzzing processes to be more efficient.

4.2 Output analyzer and Subscriber

The output analyzer takes the outputs captured by the subscriber and analyzes them to provide reports to the orchestrator.

The analyzer operates by creating a mapping of potential discontinuities and their inputs. It begins by establishing a baseline for every single output topic, by looking for integer, floating-point, Boolean, and array outputs, and stores a shortened differential vector for them. First, it collects the five (for `monte_carlo` and `spline` sampling approaches) or N (for `chebfun` sampling approach) values. Then, it starts the anomaly detection process in a separate thread while signaling to the orchestrator to send a new group of values. This ensures that the system is utilizing as much of the duty cycle as possible.

The anomaly detection algorithm is at the center of this concept. We primarily rely on numerical approximation rules, as well as analysis of previous research on bugs behavior in cyber-physical systems. Our criterium consists of: (1) numerical discontinuities *i.e.*, points where the difference between two values is significant and nonlinear, and (2) information leakage across values in a message, *i.e.*, when a field that is previously unchanged suddenly changes after input changes, especially when it is an outlier or a unique instance.

The discontinuity measures come from the underlying assumption of the continuous nature of the physical world. While there are valid reasons for discontinuity in the output of a robot, if we notice large values in the first or second derivative, mainly when such values are not reflected in their neighboring value derivatives, it necessitates further investigation. These discontinuities are often a sign of a bug in the control loop or floating-point mathematical computations. The quintessential example of this is a division that is rapidly approaching infinity as its denominator tends toward zero. We use the first derivative to look for any logical jumps, and the second derivative to look for any potential significant changes in the first derivative. The analyzer records the mean and the standard deviation for the first and second derivatives of each output. Any deviations more than a *detection threshold* are raised as a potential anomaly.

Furthermore, the analyzer includes a Bayesian sensitivity analysis[51] that calculates the weighted effect of each of the past N inputs on the output. The weighted effect is calculated by measuring derivative effects of the output from a given group of inputs. This weighted effect is stored in an array of size N, and the sum of this array is always one. The weight array is continually updating, maintaining size N. Empirically, we chose to consider the first 50% of values as

“old”, and second 50% of values as “new”. The analyzer raises an anomaly if the sum of the “new” values is less than the *weighting threshold*. This indicates the process is no longer updating that output for new inputs, flagging a potential anomaly.

We consider NaN and infinity outputs as leading special cases of discontinuity. These values are vital for analysis, given that once an internal robotic state system starts to propagate NaN values, it is highly likely that it will continuously output NaN values. Unless the system has proper handling functionality, the same applies to values of infinity. We only flag NaN and infinity values when they occur in more than three consecutive listings, as this was determined to be a sign that the robot is unable to perform any new calculations.

The analyzers also look for anomalies where changes occur unexpectedly in some outputs. Upon startup, a mapping is created for every input and output field combination. This mapping is one-to-one or one-to-many. When a new message is received from the subscriber, for every field changed a counter is incremented for those particular mappings. The analyzer raises an anomaly when the counter is incremented, and the counter value is below a *normality threshold*, a user-defined ratio of the maximum difference between any two mappings for the same input field. A low counter indicates that a value has changed when it normally remains unchanged, a potential anomaly.

By default, all three thresholds (*detection*, *weighting*, and *normality*) are set to be equal to two σ or two standard deviations from the mean. This value is based on the idea of using a 95% confidence interval for fuzzing anomaly detection from Zhao *et al.*[59]. Each threshold can be customized by the user. For example, a user can specify the *detection threshold* and *weighting threshold* at two σ , but set the *normality threshold* to three σ , if an output field has a low rate of change causing a larger number false positives than expected.

Once the anomaly detection algorithm has discovered an anomaly, the analyzer saves inputs, outputs, and timings of the test, as a pickled message object.

4.3 Crash Detection

DiscoFuzzer implements another anomaly detection mechanism, performed by the orchestrator, that fuzzers use typically: crash detection. The orchestrator deploys the ROS nodes in python subprocesses. When a node crashes, the subprocess returns with an error code that is caught by the orchestrator. Furthermore, the orchestrator periodically uses the utility *ROStopic*[38] to check whether the topics are still available to the nodes. When a crash is detected, it creates similar reports every time it discovers an anomaly.

4.4 Source Code

The source code of *DiscoFuzzer* is publicly available at URL³.

5 EVALUATION

In this section, we evaluate *DiscoFuzzer* for the effectiveness of the discontinuity-based analysis approach and efficiency of the three sampling approaches.

More specifically, we conduct an experimental campaign to answer two main research questions:

RQ1 Can *DiscoFuzzer* detect vulnerabilities in ROS nodes based on discontinuity analysis? (§5.3)

RQ1.1 How many bugs does the discontinuity analysis detect compared to the other analyses, previous work and human detection?

RQ2 Which sampling approach used by *DiscoFuzzer* is more efficient to trigger vulnerabilities in ROS nodes? (§5.4)

We initially create a set of benchmarks by including popular ROS nodes with known vulnerabilities (§5.1). Then, we execute *DiscoFuzzer* against the benchmarks (§5.2). Furthermore, we also consider some common use-cases and threats to validity to conclude the evaluation (§5.5, §5.6).

5.1 Benchmarks

Table 1 shows the benchmarks.

We identified 14 open-source ROS packages in the ROS ecosystem, using the *rosmmap*[30] tool. We chose ten of them by determining the packages that had the highest number of packages that depended on them, *i.e.*, those packages with very high impact on the community. These ten core packages are usually included in academic robotic systems to support research, development, and prototyping[24]. We chose also four additional real-world packages that are used in non-academic contexts. These packages are highly-rated on GitHub, including the Autoware self-driving car, the Udacity self-driving car, the NASA Mars rover, and the ARDrone flight system.

All code for the 14 packages is hosted and regularly updated on GitHub. For every package, we looked into its issue tracker for closed issues. We kept only those issues related to software vulnerabilities and discarded the others (*i.e.* users’ questions, feature requests, or compilation issues). We assign to each vulnerability a numeric and unique id. We define a benchmark as the code in the package’s repository at a version that include the maximum number of vulnerabilities not already included in other benchmarks. This criterion resulted in more benchmarks for the same package but at different versions. Thus, we identified 20 benchmarks with 89 vulnerabilities.

During the execution of fuzz campaign, we detected further vulnerabilities in the targets. For each of them, we extensively searched into the issue trackers and find the related issue. If not found, we consider the new vulnerability as previously unknown to the community and we submitted a new issue related to it. We assigned to all of them a numeric and unique id, and we added them to the *DiscoFuzzer*’s benchmarks. Since we reported the 8 vulnerabilities to the developers, 5 have have been officially confirmed as potential security threats. Of those 5, 3 have been patched and the other 2 are still open issues. The other 3 are still pending for confirmation.

Further information, such as the link to the issue on GitHub where the vulnerability is described, is in the JSON file at URL⁴.

³The URL is not available due to the double-blind process. The authors are available to send a copy of the file upon request.

Table 1: DiscoFuzzer’s benchmarks. The first two columns enumerate the names of the benchmarks. The third and fourth columns show the URLs of the code and the lists of the ids of each benchmark’s vulnerabilities. An exception occurs for the benchmarks number 5 (carla) and 6 (carla-bis): they require other code running to be tested, namely carla-core, introducing in the above-mentioned benchmarks two new vulnerabilities. We added the URL of carla-core in the table. We also added the previously unknown vulnerabilities in bold (cfr. §5.3).

	name	URL	vulnerabilities
1	apm-planner	https://github.com/ArduPilot/apm_planner/tree/dfe1865a82	[1, 2, 3, 4, 5, 6, 90]
2	ardupilot	https://github.com/AutonomyLab/ardrone_autonomy/tree/2e3b75a	[7, 8, 9, 10, 11, 12, 13]
3	autoware	https://github.com/autowarefoundation/autoware/tree/31f4bfb	[14, 15, 16, 17, 18, 19, 20]
4	autoware-bis	https://github.com/autowarefoundation/autoware/tree/e625625	[21]
5	carla	https://github.com/carla-simulator/ros-bridge/tree/8e468ca	[22, 23]
6	carla-bis	https://github.com/carla-simulator/ros-bridge/tree/625960e	[24]
	carla-core	https://github.com/carla-simulator/carla/tree/ec3bb90	[25, 26]
7	cartographer-ros	https://github.com/googlecartographer/cartographer_ros/tree/2538ac3	[27, 28, 29, 30, 31, 32]
8	cartographer-ros-bis	https://github.com/googlecartographer/cartographer_ros/tree/7bcdda4	[33, 34, 91]
9	cob-driver	https://github.com/ipa320/cob_driver/tree/7a5d7c8	[35, 36, 37, 38, 39, 40]
10	image-pipeline	https://github.com/ros-perception/image_pipeline/tree/d11edf3	[41, 42, 43, 44, 45, 46, 47, 48, 92]
11	lsd-slam	https://github.com/tum-vision/lsd_slam/tree/bb82258	[49, 50, 51, 52, 53, 93]
12	moveit	https://github.com/ros-planning/moveit/tree/ece11fe	[54, 55, 56, 57, 58, 59, 60, 94, 95]
13	mrpt	https://github.com/mrpt/mrpt/tree/f564006	[61, 62]
14	mrpt-bis	https://github.com/mrpt/mrpt/tree/a4bcb08	[63]
15	mrpt-tris	https://github.com/mrpt/mrpt/tree/31e853f	[64, 96]
16	navigation	https://github.com/ros-planning/navigation/tree/48323b0	[65, 66, 67, 68, 69, 70, 71, 72, 73]
17	open-source-rover	https://github.com/nasa-jpl/osr-rover-code/tree/33f072e	[74, 75, 76, 97]
18	rtabmap	https://github.com/introlab/rtabmap/tree/173bd49	[77, 78, 79, 80, 81]
19	rtabmap-bis	https://github.com/introlab/rtabmap/tree/344dc16	[82, 83, 84]
20	universal-robot	https://github.com/ros-industrial/universal_robot/tree/8c912d4	[85, 86, 87, 88, 89]

5.2 Experimental design

DiscoFuzzer’s configuration parameters were tuned with results from all three different sampling approaches. Table 2 shows the values of all other parameters that are identical among the three. These parameters were determined experimentally after exploratory analysis with DiscoFuzzer.

For each combination of a benchmark and one of the three DiscoFuzzer sampling approaches, we run the fuzzer 10 times. A single

⁴The URL is not available due to the double-blind process. The authors are available to send a copy of the file upon request.

Table 2: DiscoFuzzer’s configuration parameters used in the evaluation.

parameter	value
interval size	5
resolution	0.01
compression factor	200
detection threshold	2 times the standard deviation of the expected values
weighting threshold	3 times the standard deviation of the expected values
normality threshold	2 times the standard deviation of the expected values

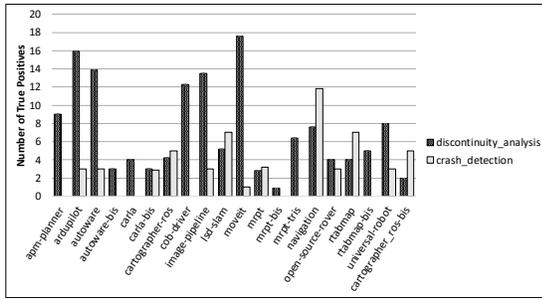
repetition, i.e., a fuzz campaign[22], lasts 24 hours. The only exception is the chebfun sampling approach: it only needs to run once per benchmark because of its deterministic nature. In total, we ran 420 fuzz campaigns, 420 full days in CPU time.

The fuzzing environment runs in a virtual machine with four cores of an Intel(R) Xeon(R) CPU E5-4650 set in a host emulation, 8 GB of RAM, Ubuntu 16.04, ROS kinetic, Python 3.6, and pychebfun⁵. We optimize the execution of the program with the python profile library and implement an automatic load balancer between the subscriber and the publisher to ensure that the system is constantly publishing at the highest rate that the fuzzed node can support. We enforce the testing time by using the system clock and terminate the program after 24 hours.

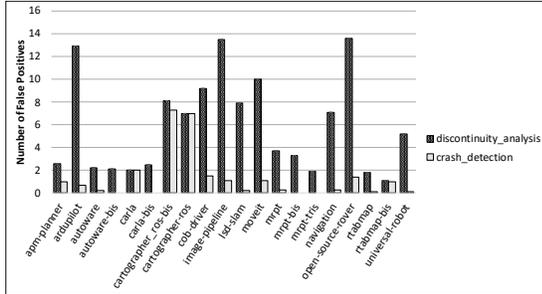
5.3 Effectiveness of DiscoFuzzer

Every fuzz campaign produced a report for each vulnerability detected by DiscoFuzzer. The report contains information such as sampling approach, target node, last messages (i/o), anomaly type, and detection time (elapsed time since the beginning of the fuzzing campaign). We manually inspected every report and tested its reproducibility. We labeled the report with its vulnerability-id if it

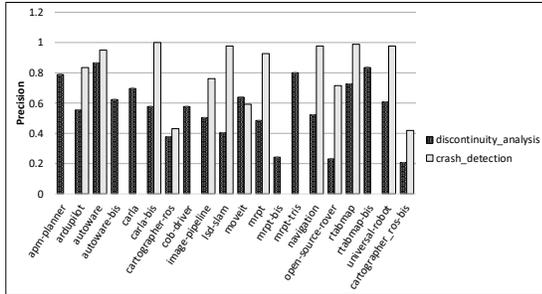
⁵The used version of pychebfun is at <https://github.com/pychebfun/pychebfun/commit/cda9283>.



(1) Number of True Positives



(2) Number of False Positives



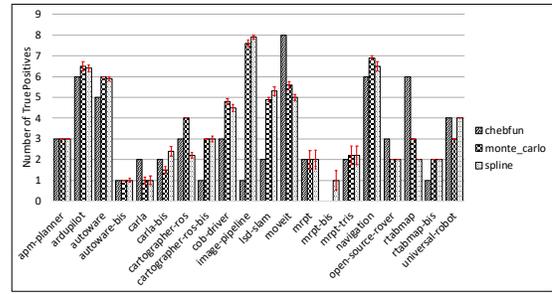
(3) DiscoFuzzer's Precision

Figure 4: Metrics grouped by the issuing detector mechanism. The bars' height represents the mean value for the group in the specified benchmark.

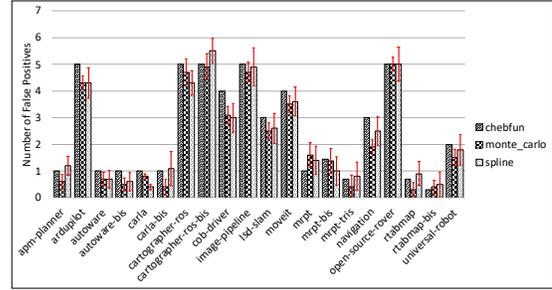
produced a True Positive. In case of a previously unknown vulnerability, we created a new one and recorded it in the benchmarks (cfr. §5.1. All of the reports consist of JSON files, uploaded as archives at URL⁶.

DiscoFuzzer found eight previously unknown vulnerabilities, marked in bold in Table 1, and identified 77 of the 89 previously known vulnerabilities. We grouped the reports by the issuing detector, i.e., either discontinuity_analysis or crash_detection. The crash detection mechanisms detected 22 distinct vulnerabilities: 20 of the 89 previously known vulnerabilities, and 2 new ones. During the evaluation, the discontinuity analysis of DiscoFuzzer detected 63 distinct vulnerabilities: 57 of the 89 previously known vulnerabilities, and 6 new ones (cfr. RQ1).

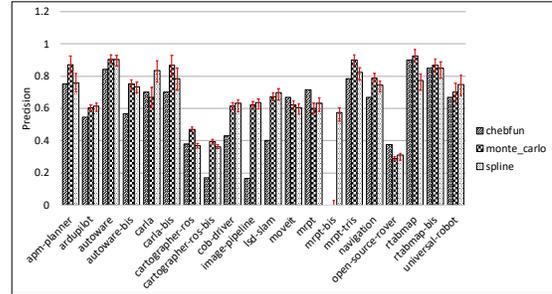
⁶The URL is not available due to the double-blind process. The authors are available to send a copy of the file upon request.



(1) Number of True Positives



(2) Number of False Positives



(3) DiscoFuzzer's Precision

Figure 5: Metrics grouped by the used sampling approach. The bars' height represents the mean value for the group in the specified benchmark. They also include an error bar, but for chebfun, that is deterministic.

Figure 4 shows the number of true and false positives, and the precision computed as the number of true positive divided by the number of all positives. Even if the discontinuity_analysis can detect more vulnerabilities in 12 out of 14 benchmarks, its complexity and infancy are highlighted by looking at the precision. The crash_detection almost has no false positives resulting in a much higher precision than discontinuity_analysis.

The discontinuity analysis of DiscoFuzzer detected 41 more unique vulnerabilities compared to the crash detection (cfr. RQ1.1).

5.4 Efficiency of the sampling approaches

We grouped the reports by the used sampling approaches, shown in Figure 5.

Following the methodology presented by Klees *et al.*[18], we performed a statistical analysis to determine which of the sampling approaches has better precision. We performed the Mann Whitney

Table 3: Results of the statistical analyses on precision performed among *DiscoFuzzer*'s three sampling approaches, namely *spline* (*s*), *monte_carlo* (*mc*), and *chebfun*(*c*). The *p* is the p-value of the Mann-Whitney U test performed, while the *A* is the Vargha and Delaney's statistic. The first column lists the benchmarks' names.

Benchmark	$p_{s,mc}$	$A_{s,mc}$	$p_{s,c}$	$A_{s,c}$	$p_{mc,c}$	$A_{mc,c}$
apm-planner	0.106	0.34	0.481	0.50	0.045	0.70
ardupilot	0.380	0.55	0.115	0.65	0.006	0.80
autoware	0.435	0.47	0.027	0.75	0.016	0.78
autoware-bis	0.432	0.47	0.099	0.66	0.049	0.70
carla	0.111	0.65	0.121	0.65	0.227	0.40
carla-bis	0.221	0.41	0.226	0.60	0.032	0.73
cartographer-ros	0.036	0.26	0.323	0.43	0.000	0.94
cartographer-ros-bis	0.154	0.36	0.000	1.00	0.000	1.00
cob-driver	0.380	0.46	0.000	1.00	0.000	1.00
image-pipeline	0.470	0.48	0.000	1.00	0.000	1.00
lsd-slam	0.351	0.56	0.000	1.00	0.000	1.00
moveit	0.296	0.42	0.057	0.30	0.035	0.30
mrpt	0.468	0.52	0.180	0.38	0.148	0.36
mrpt-bis	0.000	0.94	-	-	0.000	0.94
mrpt-tris	0.251	0.42	0.339	0.56	0.084	0.67
navigation	0.114	0.34	0.221	0.60	0.001	0.90
open-source-rover	0.350	0.45	0.220	0.40	0.000	0.00
rtabmap	0.073	0.33	0.288	0.42	0.228	0.59
rtabmap-bis	0.448	0.48	0.500	0.49	0.482	0.51
universal-robot	0.409	0.54	0.221	0.60	0.219	0.40

U test [23] to test the null hypothesis H_0 that the precision of two sampling approaches is statistically equal. In case of $p_value < 0.05$, we reject H_0 . We computed the Vargha and Delaney's A statistic [48] to compute the effect size of the statistically different group. Given two groups, $A_{1,2}$ is close to 1 if the first group has a statistically higher precision than the second group, close to 0 on the opposite case. Table 3 shows the results of this statistical analysis. The precision of the *spline* and *monte_carlo* sampling methods are statistically different only in three cases, where *spline* is more precise during a single benchmark and *monte_carlo* in other two. In 5 out of 20 benchmarks, the *spline* approach has a statistically higher precision than *chebfun*, but in one benchmark the converse is true. In 12 out of 20 benchmarks, *monte_carlo* is better than *chebfun*, but in the other two cases, the contrary holds. The *monte_carlo* sampling is the most precise while using *chebfun* leads to the least precise reports. However, in the majority of cases, the statistical tests are inconclusive.

Table 4 and Figure 6 present the vulnerabilities as detected by the three sampling approaches. Every sampling approach found unique vulnerabilities that the other two could not. In particular, the *chebfun* sampling approach found more than double the amount of these unique vulnerabilities, including two previously unknown vulnerabilities. In total, eight new vulnerabilities were detected, five of which were detected by all three methods, two which were identified by *spline* and *monte_carlo*, and two that were identified only by *chebfun*.

Table 4: Vulnerabilities as detected by the three sampling approaches. The first three columns indicate whether the row is the intersection (●) or not (○) of the elements detected by *spline*, *monte_carlo*, and *chebfun*. The fourth column has the number of vulnerabilities in the resulted subset. The last column contains the ids of the vulnerabilities, highlighting in bold the previously unknown vulnerabilities.

<i>spline</i>	<i>monte_carlo</i>	<i>chebfun</i>	#	vulnerability IDs
●	●	●	46	1, 2, 7, 8, 9, 10, 12, 13, 15, 16, 17, 21, 23, 26, 27, 34, 35, 39, 40, 53, 54, 55, 56, 57, 59, 60, 61, 62, 64, 66, 67, 68, 72, 73, 74, 81, 82, 86, 88, 89, 90, 92, 93, 96, 97
●	●	○	21	11, 14, 20, 25, 29, 32, 33, 36, 37, 41, 42, 43, 44, 45, 47, 49, 50, 52, 70, 83, 91
●	○	●	4	18, 24, 71, 85
○	●	●	3	19, 28, 80
●	○	○	3	48, 52, 63
○	●	○	2	46, 65
○	○	●	7	22, 30, 75, 77, 79, 94, 95

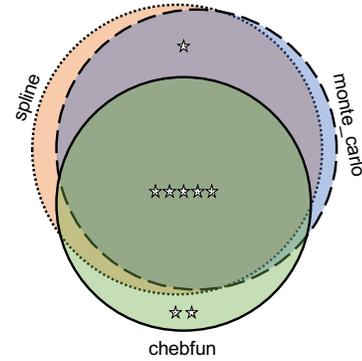


Figure 6: Venn diagram of the vulnerabilities as detected by the three sampling approaches. The area occupied by the three circles represents the number of unique vulnerabilities found. Overlapped areas indicate the number of unique vulnerabilities found by both (or the three) approaches. The stars represent unique vulnerabilities that were previously unknown.

Similarly to the precision analysis, we analyzed the vulnerability detection times of the three sampling approaches through statistical tests. The null hypothesis H_0 , for the Mann-Whitney U test, is that the detection times for a vulnerability in two sampling methods are statistically equal. The Vargha and Delaney's $A_{1,2}$ statistic is

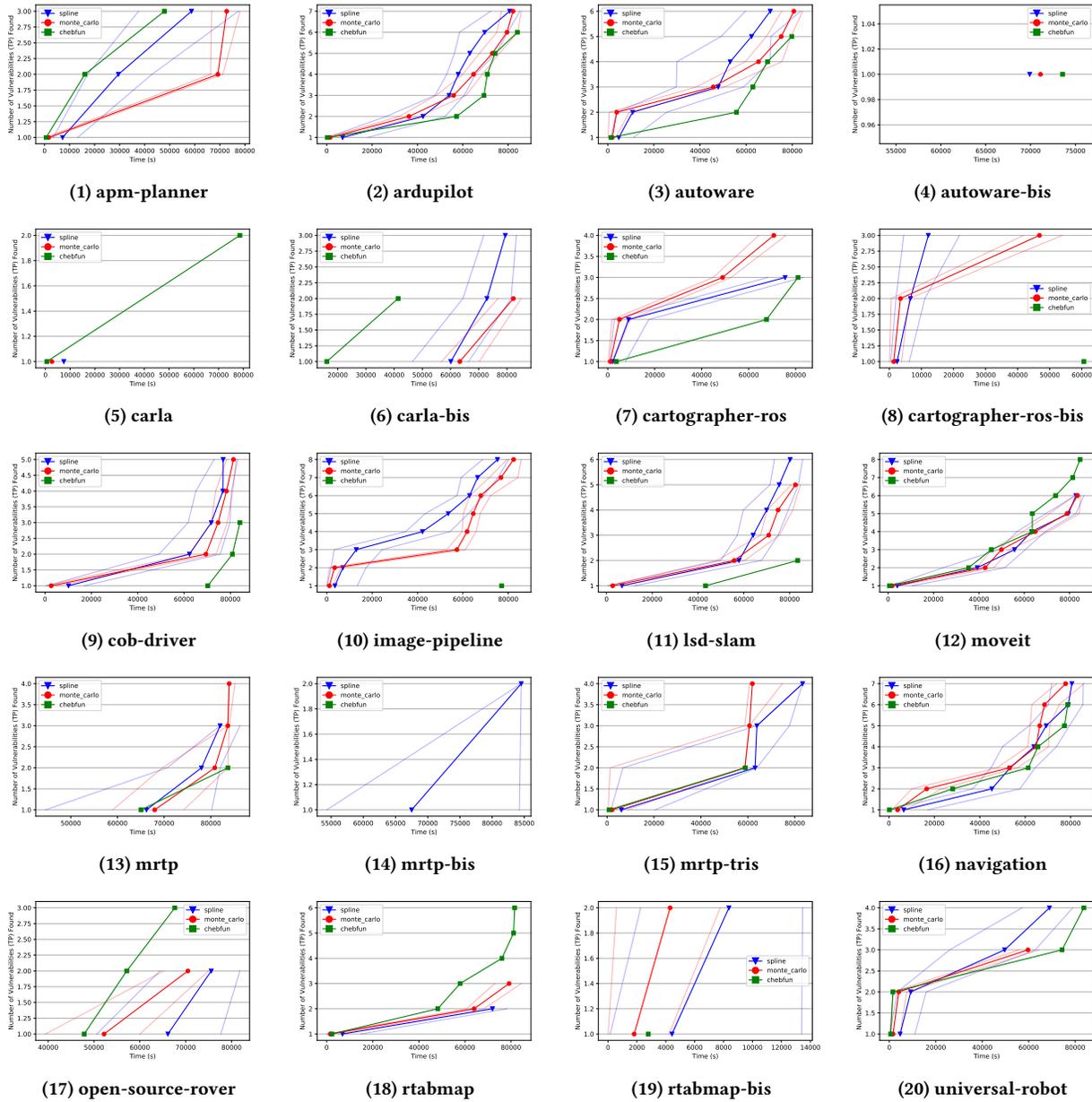


Figure 7: Detection times plotted for each benchmark (ROS node) and grouped by sampling approaches. The lighter lines around spline and monte_carlo represent the minimum and the maximum functions, while the darker lines are the average functions.

high if the first group has longer detection time than the other, low in the opposite case. We then compared each sampling approach with the other two. For each comparison, we computed how many times it is statistically faster than the other. Results of the statistical analysis on vulnerability detection times are in the appendix. Table 5 lists the results of these comparisons, demonstrating that chebfun approach is statistically faster to detect bugs. Figure 7 provides a view on the detection times for the benchmarks.

No sampling method is statistically better in finding vulnerabilities, but each of them discovered vulnerabilities the others could not. However, the chebfun approach occurs to be the fastest in finding vulnerabilities (*cf.* RQ2).

Table 5: Pairwise comparison among the three sampling approaches of *DiscoFuzzer*. While the second column present us with the results based on statistical significance, the last column presents the number of ties. There is a tie in a comparison of vulnerability detection times if such a comparison is not statistically significant, i.e., the p-value of the associated Mann-Whitney U test is greater than 0.05.

Comparison	Result	Ties
spline - monte_carlo	15 - 25	26
spline - chebfun	11 - 19	19
monte_carlo - chebfun	13 - 17	18

5.5 Discussion

DiscoFuzzer demonstrates its capabilities in detecting previously known and unknown erroneous behavior. The three sampling approaches complement each other well to identify many types of cyber-physical vulnerabilities. In the majority of cases, *DiscoFuzzer* successfully identified vulnerabilities with high precision. However, testing revealed some case-specific behaviors that are notable for future analysis.

When considering the security of ROS systems, there are three viewpoints to analyze. First of all, ROS systems are insecure and provide little protection against attackers [11]. The second point of view is that through leveraging all security modules available through ROS [52], we can treat these systems as reasonably secure. The only concern is about sensor-based attacks. As an extension of the second viewpoint, the third one only considers attacks that could potentially leave long-lasting or hard to repair damage to the physical systems or their surroundings. *DiscoFuzzer* operates under the assumption that the only vulnerabilities in ROS systems are concerning sensor-based attacks. Thus, every anomaly found by *DiscoFuzzer* is not adequately protected by the current security measures and can be exploited to cause damage to the robotic system and its surroundings. In this way, *DiscoFuzzer* provides another layer of protection.

Overall, the false positives are exceedingly rare for crash failures. Two scenarios cause them: (1) virtual memory overflow due to dynamic array sizes, and (2) asserts. Most nodes only crash due to the first scenario occurring when the fuzzer passes a message with a large array causing the virtual machine to run out of memory (as they do not have swap). Then, the fuzzer crashes but does not indicate a bug in the node code. Notably, Google Cartographer exhibits a different crash profile. Google Cartographer uses asserts while parsing messages. For any malformed message, the assert causes the node to cease execution. After a manual analysis of the design and code of the node, we hypothesize that this is a design decision made to enforce well-formatted inputs. Even applying constraints and automatically eliminating any combination that caused too many crashes, the tool still picks up a moderate amount of false positives.

In respect to false positives for discontinuity analysis, they occur primarily as a product of the threshold settings, with some behaviors producing more noise than others. However, it is necessary to note that for this paper, we chose a consistent threshold for testing

of all packages. In real-world applications, users can tune the fuzzer for their specific applications.

In general, false negatives are an expected component of any testing platform[18]. Here, we found that false-negatives are primarily due to specific configuration requirements or bugs that required multiple fields to be at precise values.

A limitation of *DiscoFuzzer* is data-type selection. *DiscoFuzzer* only analyzes numeric data types, such as integers, floating points, boolean values, and arrays. Other data types are, therefore, excluded. While our approach behaves well for our stated goals, it deliberately excludes vital data types such as strings, a typical application in fuzzing research. A future version of the tool can mitigate this limitation by using a combination approach with other well-established non-numerical fuzzers such as AFL[58]. Another limitation concerns *DiscoFuzzer*'s strict adherence to the ROS message format. Our tests never included malformed messages nor out-of-order messages. Only valid messages were sent to the node to ensure that the fuzzer was only analyzing the node-specific code and not the ROS parsing mechanisms. These concerns are ROS specific and deemed out-of-scope as the goal was to create a generalizable technique.

The final limitation is due to the black-box nature of the test system. While black-box testing allows for easy integration of a wide variety of ROS nodes written in multiple programming languages based purely on input-output behavior, an extensive corpus of literature focuses on grey-box fuzzing, which is not explored here. Grey-box fuzzing would allow for more efficient guided fuzzing, at the cost of complexity of fuzzer design and portability. While the use of grey-box fuzzing, in this case, would allow more accurate tuning of parameters, there would be a high cost in efficiency and performance. By using the black-box technique, we obtain accurate results while providing an efficient and usable tool for users.

While previous work on property-based fuzzing is valid, *DiscoFuzzer* outperforms them in detecting vulnerabilities that combine multiple properties in unison, which traditional property-based fuzzing detects less efficiently [39]. Additionally, *DiscoFuzzer* is more accessible for developers to validate discovered vulnerabilities. Similarly, another recent advancement in the field includes derivative-based fuzzing [17]. However, this approach is narrow in scope, focusing only on the control loop of systems. *DiscoFuzzer* is capable of a broader scope in its analysis, able to identify a more significant number of potential vulnerabilities with applications among any ROS package.

5.6 Threat to validity

The only threat to validity that we are aware of for *DiscoFuzzer* is external. As we focused our initial design on ROS, we did not guarantee that our discontinuity results are generalizable to other cyber-physical systems. However, we do believe that the approach is general enough to apply to any system with well-defined inputs and outputs interfaces.

6 RELATED WORK

Previous research related to our work includes fuzz testing and state-of-the-art of robot system testing.

6.1 Fuzz testing

In recent years, the security community advanced enormously in the state-of-the-art of fuzzers, proposing different approaches to generate new inputs and detect the program’s misbehavior. The most recent advancement in derivative-based fuzzing is an approach by Taegy *et al.*[17], which focused on fuzzing the PID control loop of several types of robots. This approach utilized control loop stability as a means to guide coverage testing for systems that use the MAVlink protocol. kAFL[40] is coverage-guided kernel fuzzer that is OS-independent and based on hardware-assisted instrumentation. DIFUZE[9] is an interface-aware fuzzing tool to automatically generate valid inputs and trigger the execution of the kernel drivers. JANUS[55] is a feedback-driven fuzzer that explores the two-dimensional input space of a file system, that mutates meta-data on a large image while emitting image-directed file operations. PeriScope[44] focuses on the hardware-OS boundary, targeting device drivers. Razzler[16] guides the exploration of inputs to search for data races.

Concolic execution can help the exploration of the input space where standard approaches fail. These fuzzers are called hybrid and vastly showed their efficacy in testing software programs[46], [7], [57], [59]. However, the community also provided lighter methods to explore the program under test and overcome magic numbers and nested checksums. VUzzer [33] leverages control- and data-flow features to generate interesting inputs. TFuzz[27] solves the problem by removing sanity checks in the target program when it cannot bypass them. REDQUEEN [5] exploited the input-to-state correspondence, *i.e.*, the feature of some program where parts of the input often end up directly in the program state. Some fuzzers try to use some level of knowledge of the input space and generate fewer but better inputs. TIFF [15] tags input bytes with its basic type (*e.g.*, 32-bit integer) in the program, then it uses this information to mutate the inputs accordingly. ProFuzzer [56] adds to this the ability to understand input fields of critical importance. Superior [50] creates meaningful tests for programs that process structured input (*e.g.*, XML engines) by analyzing its grammar.

The differential testing approach is a way to find errors by comparing results from different implementations of the same program, resulting in cross-referencing oracles. Chen *et al.*[8] applied differential testing to a coverage-guided fuzzer for Java Virtual Machines. Nezha [29] exploits these differences also to generate inputs that are more likely to trigger a bug. DiffFuzz [26] explores the input space by maximizing the differences in resource consumption since it targets side-channels vulnerabilities.

Sanitizers are another approach to detect specific failures in a fuzzing campaign. AddressSanitizer[41] finds out-of-bounds accesses to heap, stack, and global objects, as well as use-after-free bugs. UndefinedBehaviorSanitizer [21] detects null pointer, signed integer overflows, and wrong type conversions. MemorySanitizer [20] focuses on uninitialized reads.

The novelty of *DiscoFuzzer* lies in taking the previous model-based testing approaches and extending their most commonly found failure modes into a highly search-based system. We provide a standard extension framework for analyzing model-based bugs. *DiscoFuzzer* is highly user-customizable and aimed at an under-explored area of fuzzer research: cyber-physical systems. As of

the publishing of this paper, this is the first-known approach that combines these methods for applications within the cyber-physical domain. Additionally, this approach is the first applied to Robotic systems, namely ROS.

6.2 Robotic System Testing

Robotic systems are validated with formal specifications, and these formal verification approaches are complementary to the testing of the actual implementation of the robot. With the assumption of functioning hardware, we focus on how to test robotic software.

DiscoFuzzer falls into the category of search-based approaches. A search-based method uses heuristics to explore the input space of the target and automatically generate test cases. Both Ali and Yule[2] and Arrieta *et al.*[3][4] use genetic algorithms to generate and select test cases. Matinnejad *et al.*[25] applied different search algorithms for testing automotive embedded systems, including random search, adaptive random search, a hill-climbing algorithm, and a simulated annealing algorithm. Abbas *et al.*[1] applies the Monte Carlo simulation to explore the input space to find violations of robustness properties. Instead, *DiscoFuzzer* is not applying Monte Carlo on models nor focuses on temporal logic. *DiscoFuzzer* treats the module as a black-box, and its Monte Carlo implementation is aware of the state of the robot system.

Specific to ROS, RosPenTo is a semi-automated tool for testing ROS[11]. It injects fake messages into the system to demonstrate the consequences of the lack of security in ROS. Indeed, an attacker (1) can easily query the master for sensitive information, and (2) can easily impersonate a subscriber node. These are the assumptions we used for *DiscoFuzzer*. Indeed, *DiscoFuzzer* assumes that the message and packaging interface of ROS are wholly intact and reliable, and focuses on the effects of those messages on the system.

Santos *et al.*[39] implemented a property-based testing framework for ROS. This first approach aims at the automatic generation of test scripts for property-based testing of various configurations of a ROS system. Their approach looks for sequences of messages that either crashes the target node or violate a previously specified property. An automatic test generation method builds the property-based tests from configuration models extracted by a static analysis framework. Instead, *DiscoFuzzer* is a proper fuzzer that, besides automatically generated test cases for the target, does not need any specification of properties to detect anomalies in the target.

7 CONCLUSION

In this paper, we have proposed a novel fuzz testing methodology. We chose three different sampling methods and implemented the methodology in *DiscoFuzzer* targeting ROS packages. We tested it against 14 ROS packages to evaluate its efficacy and efficiency. Results show that *DiscoFuzzer* can find more vulnerabilities than traditional crash detection mechanisms without the necessity to specify any ad-hoc property for the targets. However, no sampling approach resulted entirely better than the other, but all of them are necessary to find the majority of unique vulnerabilities. Thus, we envision *DiscoFuzzer* will be able to use a combined approach for sampling and cover all the detected vulnerabilities at once. Furthermore, future research should look forward to improving discontinuity-based analysis and decrease the number of false positives while using the approach.

REFERENCES

- [1] Houssam Abbas, Georgios Fainekos, Sriram Sankaranarayanan, Franjo Ivančić, and Aarti Gupta. 2013. Probabilistic temporal logic falsification of cyber-physical systems. *ACM Transactions on Embedded Computing Systems (TECS)* 12, 2s (2013), 1–30.
- [2] Shaukat Ali and Tao Yue. 2015. U-test: evolving, modelling and testing realistic uncertain behaviours of cyber-physical systems. In *2015 IEEE 8th International Conference on Software Testing, Verification and Validation (ICST)*. IEEE, 1–2.
- [3] Aitor Arrieta, Shuai Wang, Gouriya Sagardui, and Leire Etxeberria. 2016. Search-based test case selection of cyber-physical system product lines for simulation-based validation. In *Proceedings of the 20th International Systems and Software Product Line Conference*. 297–306.
- [4] Aitor Arrieta, Shuai Wang, Gouriya Sagardui, and Leire Etxeberria. 2016. Test case prioritization of configurable cyber-physical systems with weight-based search algorithms. In *Proceedings of the Genetic and Evolutionary Computation Conference 2016*. 1053–1060.
- [5] Cornelius Aschermann, Sergej Schumilo, Tim Blazytko, Robert Gawlik, and Thorsten Holz. 2019. REDQUEEN: Fuzzing with Input-to-State Correspondence. In *NDSS*, Vol. 19. 1–15.
- [6] Mordechai Ben-Ari and Francesco Mondada. 2018. Robots and their applications. In *Elements of Robotics*. Springer, 1–20.
- [7] Peng Chen and Hao Chen. 2018. Angora: Efficient fuzzing by principled search. In *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE, 711–725.
- [8] Yuting Chen, Ting Su, Chengnian Sun, Zhendong Su, and Jianjun Zhao. 2016. Coverage-directed differential testing of JVM implementations. In *proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation*. 85–99.
- [9] Jake Corina, Aravind Machiry, Christopher Salls, Yan Shoshitaishvili, Shuang Hao, Christopher Kruegel, and Giovanni Vigna. 2017. Difuze: Interface aware fuzzing for kernel drivers. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. 2123–2138.
- [10] Drew Davidson, Hao Wu, Rob Jellinek, Vikas Singh, and Thomas Ristenpart. 2016. Controlling UAVs with sensor input spoofing attacks. In *10th {USENIX} Workshop on Offensive Technologies ({WOOT} 16)*.
- [11] Bernhard Dieber, Benjamin Breiling, Sebastian Taurer, Severin Kacianka, Stefan Rass, and Peter Scharfner. 2017. Security for the robot operating system. *Robotics and Autonomous Systems* 98 (2017), 192–203.
- [12] Ayssam Elkady and Tarek Sobh. 2012. Robotics middleware: A comprehensive literature survey and attribute-based bibliography. *Journal of Robotics* 2012 (2012).
- [13] MA Hernández. 2001. Chebyshev’s approximation algorithms and applications. *Computers & Mathematics with Applications* 41, 3-4 (2001), 433–445.
- [14] ISO 8373:2012(en) 2012. *Robots and robotic devices – Vocabulary*. Standard. International Organization for Standardization, Geneva, CH.
- [15] Vivek Jain, Sanjay Rawat, Cristiano Giuffrida, and Herbert Bos. 2018. TIFF: using input type inference to improve fuzzing. In *Proceedings of the 34th Annual Computer Security Applications Conference*. 505–517.
- [16] Dae R Jeong, Kyungtae Kim, Basavesh Shivakumar, Byoungyoung Lee, and Insik Shin. 2019. Razzler: Finding kernel race bugs through fuzzing. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 754–768.
- [17] Taegy Kim, Chung Hwan Kim, Junghwan Rhee, Fan Fei, Zhan Tu, Gregory Walkup, Xiangyu Zhang, Xinyan Deng, and Dongyan Xu. 2019. RVFUZZER: finding input validation bugs in robotic vehicles through control-guided testing. In *28th {USENIX} Security Symposium ({USENIX} Security 19)*. 425–442.
- [18] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. 2018. Evaluating fuzz testing. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 2123–2138.
- [19] Sofiane Lagraa, Maxime Cailac, Sean Rivera, Frédéric Beck, and Radu State. 2019. Real-time attack detection on robot cameras: A self-driving car application. In *2019 Third IEEE International Conference on Robotic Computing (IRC)*. IEEE, 102–109.
- [20] LLVM-Developers. 2020. MemorySanitizer. (2020). <http://clang.llvm.org/docs/MemorySanitizer.html>
- [21] LLVM-Developers. 2020. UndefinedBehaviorSanitizer. (2020). <https://clang.llvm.org/docs/UndefinedBehaviorSanitizer.html>
- [22] Valentin Jean Marie Manès, HyungSeok Han, Choongwoo Han, Sang Kil Cha, Manuel Egele, Edward J Schwartz, and Maverick Woo. 2019. The art, science, and engineering of fuzzing: A survey. *IEEE Transactions on Software Engineering* (2019).
- [23] Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics* (1947), 50–60.
- [24] Aaron Martinez and Enrique Fernández. 2013. *Learning ROS for robotics programming*. Packt Publishing Ltd.
- [25] Reza Matinnejad, Shiva Nejati, Lionel Briand, Thomas Bruckmann, and Claude Poull. 2015. Search-based automated testing of continuous controllers: Framework, tool support, and case studies. *Information and Software Technology* 57 (2015), 705–722.
- [26] Shirin Nilizadeh, Yannic Noller, and Corina S Pasareanu. 2019. DiffFuzz: differential fuzzing for side-channel analysis. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 176–187.
- [27] Hui Peng, Yan Shoshitaishvili, and Mathias Payer. 2018. T-Fuzz: fuzzing by program transformation. In *2018 IEEE Symposium on Security and Privacy (SP)*. IEEE, 697–710.
- [28] Deborah Petrarra. 2019. The Rise of ROS: Nearly 55% of total commercial robots shipped in 2024 Will Have at Least One Robot Operating System package. (2019). <https://www.bloomberg.com/press-releases/2019-05-16/the-rise-of-ros-nearly-55-of-total-commercial-robots-shipped-in-2024-will-have-at-least-one-robot-operating-system-package>
- [29] Theofilos Petsios, Adrian Tang, Salvatore Stolfo, Angelos D Keromytis, and Suman Jana. 2017. Nezha: Efficient domain-independent differential testing. In *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 615–632.
- [30] Marc Pichler, Bernhard Dieber, and Martin Pinzger. 2019. Can i depend on you? Mapping the dependency and quality landscape or ROS packages. In *Proceedings of the 3rd International Conference on Robotic Computing*. IEEE.
- [31] Morgan Quigley, Ken Conley, Brian Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, and Andrew Y Ng. 2009. ROS: an open-source Robot Operating System. In *ICRA workshop on open source software*, Vol. 3. Kobe, Japan, 5.
- [32] Ragunathan Rajkumar, Insup Lee, Lui Sha, and John Stankovic. 2010. Cyber-physical systems: the next computing revolution. In *Design Automation Conference*. IEEE, 731–736.
- [33] Sanjay Rawat, Vivek Jain, Ashish Kumar, Lucian Cojocar, Cristiano Giuffrida, and Herbert Bos. 2017. VUzzer: Application-aware Evolutionary Fuzzing. In *NDSS*, Vol. 17. 1–14.
- [34] Sean Rivera, Sofiane Lagraa, Antonio Ken Iannillo, and Radu State. 2019. Auto-encoding Robot State against Sensor Spoofing Attacks. In *2019 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)*. IEEE, 252–257.
- [35] ROS-Industrial. 2020. ROS-Industrial. (2020). <https://rosindustrial.org/>
- [36] ROS.org. 2020. Automatic Testing with ROS. (2020). <http://wiki.ros.org/Quality/Tutorials/UnitTesting>
- [37] ROS.org. 2020. ROS running on ISS. (2020). <https://www.ros.org/news/2014/09/ros-running-on-iss.html>
- [38] ROS.org. 2020. rostopic. <http://wiki.ros.org/rostopic>. (2020).
- [39] André Santos, Alcino Cunha, and Nuno Macedo. 2018. Property-based testing for the robot operating system. In *Proceedings of the 9th ACM SIGSOFT International Workshop on Automating TEST Case Design, Selection, and Evaluation*. 56–62.
- [40] Sergej Schumilo, Cornelius Aschermann, Robert Gawlik, Sebastian Schinzel, and Thorsten Holz. 2017. kaf: Hardware-assisted feedback fuzzing for {OS} kernels. In *26th {USENIX} Security Symposium ({USENIX} Security 17)*. 167–182.
- [41] Konstantin Serebryany, Derek Bruening, Alexander Potapenko, and Dmitriy Vyukov. 2012. AddressSanitizer: A fast address sanity checker. In *Presented as part of the 2012 {USENIX} Annual Technical Conference ({USENIX} {ATC} 12)*. 309–318.
- [42] Amazon Web Services. 2020. Announcing AWS RoboMaker: A New Cloud Robotics Service. (2020). <https://aws.amazon.com/about-aws/whats-new/2018/11/announcing-aws-robomaker-a-new-cloud-robotics-service/>
- [43] Alexander Shapiro. 2003. Monte Carlo sampling methods. *Handbooks in operations research and management science* 10 (2003), 353–425.
- [44] Dokyung Song, Felicitas Hetzelt, Dipanjan Das, Chad Spensky, Yeoul Na, Stijn Volckaert, Giovanni Vigna, Christopher Kruegel, Jean-Pierre Seifert, and Michael Franz. 2019. PeriScope: An Effective Probing and Fuzzing Framework for the Hardware-OS Boundary. In *NDSS*.
- [45] Gerald Steinbauer. 2012. A survey about faults of robots used in robocup. In *Robot Soccer World Cup*. Springer, 344–355.
- [46] Nick Stephens, John Grosen, Christopher Salls, Andrew Dutcher, Ruoyu Wang, Jacopo Corbetta, Yan Shoshitaishvili, Christopher Kruegel, and Giovanni Vigna. 2016. Driller: Augmenting Fuzzing Through Selective Symbolic Execution. In *NDSS*, Vol. 16. 1–16.
- [47] Chris Swierczewski and Olivier Verdier. 2020. Pychebfun. <https://github.com/pychebfun/pychebfun>. (2020).
- [48] András Vargha and Harold D Delaney. 2000. A critique and improvement of the CL common language effect size statistics of McGraw and Wong. *Journal of Educational and Behavioral Statistics* 25, 2 (2000), 101–132.
- [49] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, CJ Carey, Ilhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R Harris, Anne M. Archibald, António H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods* (2020). <https://doi.org/10.1038/s41592-019-0686-2>

- [50] Junjie Wang, Bihuan Chen, Lei Wei, and Yang Liu. 2019. Superior: Grammar-aware greybox fuzzing. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 724–735.
- [51] Robert Weiss. 1996. An approach to Bayesian sensitivity analysis. *Journal of the Royal Statistical Society: Series B (Methodological)* 58, 4 (1996), 739–750.
- [52] Ruffin White, Dr Christensen, I Henrik, Dr Quigley, et al. 2016. SROS: Securing ROS over the wire, in the graph, and through the kernel. *arXiv preprint arXiv:1611.07060* (2016).
- [53] Johannes Wienke and Sebastian Wrede. 2017. Results of the survey: failures in robotics and intelligent systems. *arXiv preprint arXiv:1708.07379* (2017).
- [54] Windows. 2020. Windows 10 IoT + ROS. (2020). <https://microsoft.github.io/Win-RoS-Landing-Page/>
- [55] Wen Xu, Hyungon Moon, Sanidhya Kashyap, Po-Ning Tseng, and Taesoo Kim. 2019. Fuzzing file systems via two-dimensional input space exploration. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 818–834.
- [56] Wei You, Xueqiang Wang, Shiqing Ma, Jianjun Huang, Xiangyu Zhang, XiaoFeng Wang, and Bin Liang. 2019. Profuzzer: On-the-fly input type probing for better zero-day vulnerability discovery. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 769–786.
- [57] Insu Yun, Sangho Lee, Meng Xu, Yeongjin Jang, and Taesoo Kim. 2018. {QSYM}: A practical concolic execution engine tailored for hybrid fuzzing. In *27th {USENIX} Security Symposium ({USENIX} Security 18)*. 745–761.
- [58] Michal Zalewski. 2015. american fuzzy lop technical “whitepaper”. URL: http://lcamtuf.coredump.cx/afl/technical_details.txt (2015).
- [59] Lei Zhao, Yue Duan, Heng Yin, and Jifeng Xuan. 2019. Send Hardest Problems My Way: Probabilistic Path Prioritization for Hybrid Fuzzing.. In *NDSS*.